

1. (30 points) *PCA via successive deflation*. Let  $\theta_1, \dots, \theta_k$  be the eigenvectors of the  $k$  largest eigenvalues of  $\mathbf{C} = (1/n)\Phi\Phi^\top$ , i.e., the PCA solution. These satisfy

$$\theta_c^\top \theta_d = \begin{cases} 0 & c \neq d \\ 1 & c = d. \end{cases}$$

We will construct a method of finding  $\theta_c$  sequentially.

As we showed in class,  $\theta_1$  is the first principal eigenvector of  $\mathbf{C}$ , and satisfies  $\mathbf{C}\theta_1 = \lambda_1\theta_1$ . Now define  $\tilde{\phi}_i$  as the orthogonal projection of  $\phi_i$  onto the space orthogonal to  $\theta_1$ :

$$\tilde{\phi}_i = (\mathbf{I} - \theta_1\theta_1^\top)\phi_i.$$

Define  $\tilde{\Phi} = [\phi_1 \ \cdots \ \phi_n]$  as the *deflated* data matrix, we have

$$\tilde{\Phi} = (\mathbf{I} - \theta_1\theta_1^\top)\Phi.$$

- (a) Using the fact that  $\mathbf{C}\theta_1 = \lambda_1\theta_1$  and  $\|\theta_1\|^2 = 1$ , show that the second moment of the deflated matrix is given by

$$\tilde{\mathbf{C}} = \frac{1}{n}\tilde{\Phi}\tilde{\Phi}^\top = \mathbf{C} - \lambda_1\theta_1\theta_1^\top.$$

- (b) Show that  $\theta_2$  is the eigenvector of the largest eigenvalue of  $\tilde{\mathbf{C}}$ . Recall that  $\theta_2$  is the eigenvector of the second-largest eigenvalue of  $\mathbf{C}$ .
- (c) Suppose we have a simple method for finding the leading eigenvector and eigenvalue of a positive semi-definite matrix, denoted by  $(\lambda, \mathbf{u}) = f(\mathbf{C})$ . Write some pseudo code for finding the first  $k$  principal basis vectors of  $\Phi$  that only uses the special  $f$  function and simple vector arithmetic, i.e., your code should not use SVD or the eig function. Hint: this should be a simple iterative routine that takes 2–3 lines to write. The input is  $\mathbf{C}$ ,  $k$  and the function  $f$ , the output should be  $\theta_c$  and  $\lambda_c$  for  $c = 1, \dots, k$ . Do not worry about being syntactically correct.

Ans. 1  $\Rightarrow$

(a) = given that  $\Rightarrow$

$$\theta_c^\top \theta_d = \begin{cases} 0 & c \neq d \\ 1 & c = d \end{cases}$$

So we know  $\boxed{\theta_1^\top \theta_1 = 1}$  ( $c=d=1$ )

①

also given that  $\Rightarrow$

$$\boxed{\tilde{\phi} = (I - \theta_1 \theta_1^T) \phi} \rightarrow \textcircled{2}$$

from matrix properties we know that  $\Rightarrow$

$$(\gamma\gamma)^T = \gamma^T \gamma^T$$

and  $(\lambda)^T = \lambda$  (if  $\lambda$  is symmetric)

Here we know that  $\Rightarrow$

$$(I - \theta_1 \theta_1^T) = \text{Symmetric matrix}$$

Hence

$$\boxed{(I - \theta_1 \theta_1^T)^T = (I - \theta_1 \theta_1^T)} \rightarrow \textcircled{3}$$

To prove  $\Rightarrow$

$$\tilde{C} = \frac{1}{n} \tilde{\phi} \tilde{\phi}^T = C - \alpha_1 \theta_1 \theta_1^T$$

so,

$$\tilde{C} = \frac{1}{n} \tilde{\phi} \tilde{\phi}^T$$

Putting value of  $\tilde{\phi}$  from equation ②  $\Rightarrow$

$$\tilde{C} = \frac{1}{n} \left[ (\mathbb{I} - \theta_1 \theta_1^T) \phi ((\mathbb{I} - \theta_1 \theta_1^T) \phi)^T \right]$$

using fact  $(xy)^T = y^T x^T$ ,  $(\mathbb{I} - \theta_1 \theta_1^T)$  is symmetric

$$\tilde{C} = \frac{1}{n} \left[ (\mathbb{I} - \theta_1 \theta_1^T) \phi \phi^T (\mathbb{I} - \theta_1 \theta_1^T) \right]$$

$$\tilde{C} = \frac{1}{n} \left[ (\mathbb{I} \cdot \phi \phi^T - \theta_1 \theta_1^T \phi \phi^T) (\mathbb{I} - \theta_1 \theta_1^T) \right]$$

$$\tilde{C} = \frac{1}{n} \left[ (\phi \phi^T - \theta_1 \theta_1^T \phi \phi^T - \phi \phi^T \theta_1 \theta_1^T + \theta_1 \theta_1^T \phi \phi^T \theta_1 \theta_1^T) \right]$$

given that  $c\theta_1 = \lambda_1\theta_1$

$$\left(\frac{1}{n}\phi\phi^T\right)\theta_1 = \lambda_1\theta_1$$

$$\boxed{\phi\phi^T\theta_1 = n\lambda_1\theta_1} \rightarrow \textcircled{4}$$

and  $(\phi\phi^T\theta_1)^T = (n\lambda_1\theta_1)^T$

$$\boxed{\theta_1^T\phi\phi^T = n\lambda_1\theta_1^T} \rightarrow \textcircled{5}$$

Putting values from \textcircled{4} and \textcircled{5} we get =

$$\tilde{C} = \frac{1}{n} [\phi\phi^T - \theta_1 n\lambda_1\theta_1^T - n\lambda_1\theta_1\theta_1^T + \theta_1 n\lambda_1\theta_1^T\theta_1\theta_1^T]$$

since  $[\theta_1^T\theta_1 = 1 \rightarrow \text{from equation } \textcircled{1}]$

$$\tilde{C} = \frac{1}{n} [\phi\phi^T - \theta_1 n\lambda_1\theta_1^T - n\lambda_1\theta_1\theta_1^T + \theta_1 n\lambda_1 \cdot 1 \cdot \theta_1^T]$$

$$\tilde{C} = \frac{1}{n} [\phi\phi^T - \cancel{\theta_1 n\lambda_1\theta_1^T} - n\lambda_1\theta_1\theta_1^T + \cancel{\theta_1 n\lambda_1\theta_1^T}]$$

$$\tilde{C} = \frac{1}{n} [\phi\phi^T - n\lambda_1\theta_1\theta_1^T]$$

$$\tilde{C} = \frac{1}{n} \phi\phi^T - \lambda_1\theta_1\theta_1^T$$

$$\tilde{C} = C - \lambda_1\theta_1\theta_1^T$$

Hence  
proved

Hence we can say that =

$$\tilde{C} = \frac{1}{n} \tilde{\phi}\tilde{\phi}^T = C - \lambda_1\theta_1\theta_1^T$$

(b) = we have proved in part ① =

$$\tilde{C} = \left( \frac{1}{n} \phi \phi^T - \lambda_1 \theta_1 \theta_1^T \right)$$

multiplying both sides by  $\theta_j \neq$

$$\tilde{C} \theta_j = \left( \frac{1}{n} \phi \phi^T - \lambda_1 \theta_1 \theta_1^T \right) \theta_j$$

$$\tilde{C} \theta_j = \frac{1}{n} \phi \phi^T \theta_j - \lambda_1 \theta_1 \theta_1^T \theta_j$$

from part ① we know  $(\phi \phi^T \theta_j = n \lambda_1 \theta_j)$

so we get,

$$\tilde{C} \theta_j = \lambda_j \theta_j - \lambda_1 \theta_1 \theta_1^T \theta_j$$

also  $[\theta_1^T \theta_j = 0, \text{ hence } \lambda_1 \theta_1 \theta_1^T \theta_j = 0]$   
if ( $j \neq 1$ )

Hence,

(when  $j \neq 1$ )

$$\boxed{\tilde{C} \theta_j = \lambda_j \theta_j}$$

when ( $j=1$ )

$$\tilde{C}\theta_1 = \lambda_1\theta_1 - \lambda_1(\theta_1\theta_1^\top\theta_1) \quad [\theta_1^\top\theta_1 = 1]$$

$$\tilde{C}\theta_1 = \lambda_1\theta_1 - \lambda_1\theta_1$$

$$\tilde{C}\theta_1 = (\lambda_1 - \lambda_1)\theta_1$$

$$\tilde{C}\theta_1 = 0\theta_1$$

$$\boxed{\tilde{C}\theta_1 = 0}$$

Hence for ( $j \neq 1$ )  $\theta_j$  is a principle eigenvector of  $\tilde{C}$  with same eigenvalue  $\lambda_j$ . Also  $\theta_1$  is an eigenvector of  $\tilde{C}$  with eigenvalue 0.

Since  $\theta_1, \theta_2, \dots, \theta_k$  are the first  $k$  eigenvectors with largest eigenvalues of  $C$  i.e the principal basis vectors therefore

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$$

We know that for  $\tilde{C}$ ,  $\theta_j$  are the principal eigenvectors with eigenvalues  $(\theta_1, \lambda_1, \lambda_2, \dots, \lambda_k)$ . Therefore from above equation we can say that  $\lambda_2$  is the largest eigenvalue of  $\tilde{C}$  (since  $\lambda_1$  is not an eigenvalue of  $\tilde{C}$ ).

Hence  $\theta_2$  is the first principle eigenvector.

(C) = Pseudocode to find first k principal  
basis vectors of  $\phi$ .

function obtain Principal Basis Vector ( $C_1 K_1 f$ ):

$l_{int} - \lambda = []$

$l_{int} - u = []$

    for loop i in range (k):

$\lambda, u = f(C)$

$C = C - \lambda * u * u^T$  transpose

$(l_{int} - \lambda). append (\lambda)$

$(l_{int} - u). append (u)$

    return  $(l_{int} - u), (l_{int} - \lambda)$

2. (20 points) Consider a Gaussian mixture model in which the marginal distribution of the latent variable  $\mathbf{y}$  is  $\Pr(\mathbf{y} = \mathbf{e}_c) = \pi_c, c = 1, \dots, k$ , and the conditional distribution of  $\mathbf{x}$  given  $\mathbf{y}$  is  $\mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma})$ , i.e., each Gaussian component has their own mean  $\boldsymbol{\mu}_c$  but they share the same covariance matrix  $\boldsymbol{\Sigma}$ . Given a set of observations  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , derive the expectation-maximization algorithm for estimating the model parameters  $\pi_c, \boldsymbol{\mu}_c, c = 1, \dots, k$ , and  $\boldsymbol{\Sigma}$ .

Anw.2 =

## Expectation - Maximization Algorithm =

The Expectation - maximization algorithm is an iterative method to find the M.L.E- or M.A.P. estimate for models with latent variables.

Stepn for EM algonithm =

i) = Initialization = first all randomly initialize parameters  
 $\theta^0$  (all the  $\boldsymbol{\mu}_c$ ,  $\boldsymbol{\Sigma}$  and  $\pi_c$  values)

$\&$  = Initialize ( $\boldsymbol{\mu}$ ,  $\sigma$ ,  $\pi$ ) with random values.

$$\boldsymbol{\mu} = [\text{rand}()] * K$$

$$\sigma = [\text{rand}()] * K$$

$$\pi = [\text{rand}()] * K$$

$K =$   
 number  
 of  
 gaussians

2)- Expectation Step = By Assuming that parameter ( $\theta^{t-1}$ ) value from previous step are fixed we compute the expected values of the latent variable. (or a function of the expected values of latent variables)

calculating Expectation ( $\Psi_{ic}$ )  $\Rightarrow$

$$\Psi_{ic} = \frac{\pi_c N(x_i | \mu_c, \Sigma)}{\sum_{j=1}^k \pi_j N(x_i | \mu_j, \Sigma)}$$

where  $x_i$  = observations.

$(\Sigma, \mu, \pi)$  = parameters value from step ①.

(3) = Maximization Step - from the given values we computed in last step (values of latent variables) we estimate new values ( $\Theta^+$ ) that maximize a . variant of likelihood function.

Re-estimate and update parameters using current  $\Psi_{ic}$  from previous step.

$$\mu_c^{\text{new}} = \frac{1}{\sum_{i=1}^n \Psi_{ic}} \left[ \sum_{i=1}^n \Psi_{ic} x_i \right]$$

$$\pi_c^{\text{new}} = \frac{\sum_{i=1}^n \Psi_{ic}}{n}$$

$$\Sigma^{\text{new}} = \frac{\sum_{i=1}^n \sum_{c=1}^K \Psi_{ic} (x_i - \mu_c) (x_i - \mu_c)^T}{\sum_{i=1}^n \sum_{c=1}^K \Psi_{ic}}$$

④ = Exit condition = if likelihood of observations have not changed much then we stop the iterations otherwise we go back to Step ② and Repeat steps ② and ③ until convergence.

Suppose  $L(\alpha, \mu_C, \Sigma, \tau_C)$  is a function to compute log-likelihoods.

Then we will calculate log likelihood for current parameters and previous parameters and compare them against a threshold (tolerated change for log likelihood)

$\text{old\_L} = L(\alpha, \mu_C, \Sigma, \tau_C)$ $\text{new\_L} = L(\alpha, \mu_C^{\text{new}}, \Sigma^{\text{new}}, \tau_C^{\text{new}})$ if absolute ( $\text{old\_L} - \text{new\_L}$ ) < tolerated threshold then Exit else go back to Step ②.	$(\text{Ex} = 0.001)$
---	-----------------------

3. (50 points) *20 Newsgroup revisited.* Let us revisit the 20 Newsgroup data set <<http://qwone.com/~jason/20Newsgroups/>>, and apply some of the unsupervised methods by ignoring their labels. We will only consider the training data. You are required to code the algorithms by yourselves in the language of your choice.
- (a) *LSI/PCA via orthogonal iteration.* Implement the orthogonal iteration algorithm that finds the PCA projection matrix  $\Theta$  of a data matrix  $\Phi$ . You are allowed to use a pre-existing function of QR. Apply tf-idf to the term-document matrix to obtain  $\Phi$  and feed it into your orthogonal iteration algorithm. Remember to use sparse matrix operations to avoid unnecessary memory/computational complexities. Set  $k = 2$  and let the algorithm run until  $\Theta$  doesn't change much. Then get  $\mathbf{Y} = \Theta^\top \Phi$ . Each column of  $\mathbf{Y}$  is a two-dimensional vector that you can plot on a plain. Plot all the documents on a two-dimensional plain, and use a different color for each point that belong to different news groups.
- (b) *GMM via EM.* Implement the EM algorithm for the Gaussian mixture model (with different means and covariances for each Gaussian component). The data matrix  $\Phi$  is obtained from LSI with  $k_{\text{LSI}} = 100$  using the previous orthogonal iteration algorithm. Run the EM algorithmn for GMM with  $k_{\text{GMM}} = 20$  until convergence. For each Gaussian component  $\mathcal{N}(\mu_c, \Sigma_c)$ , calculate  $\Theta\mu_c$  where  $\Theta$  is the PCA projection; the vector  $\Theta\mu_c$  should be element-wise nonnegative. For each cluster  $c$ , show the 10 terms that have the highest value in  $\Theta\mu_c$ . The index-term mapping can be found here <<http://qwone.com/~jason/20Newsgroups/vocabulary.txt>>. Does the result make sense?

Ans.3 =

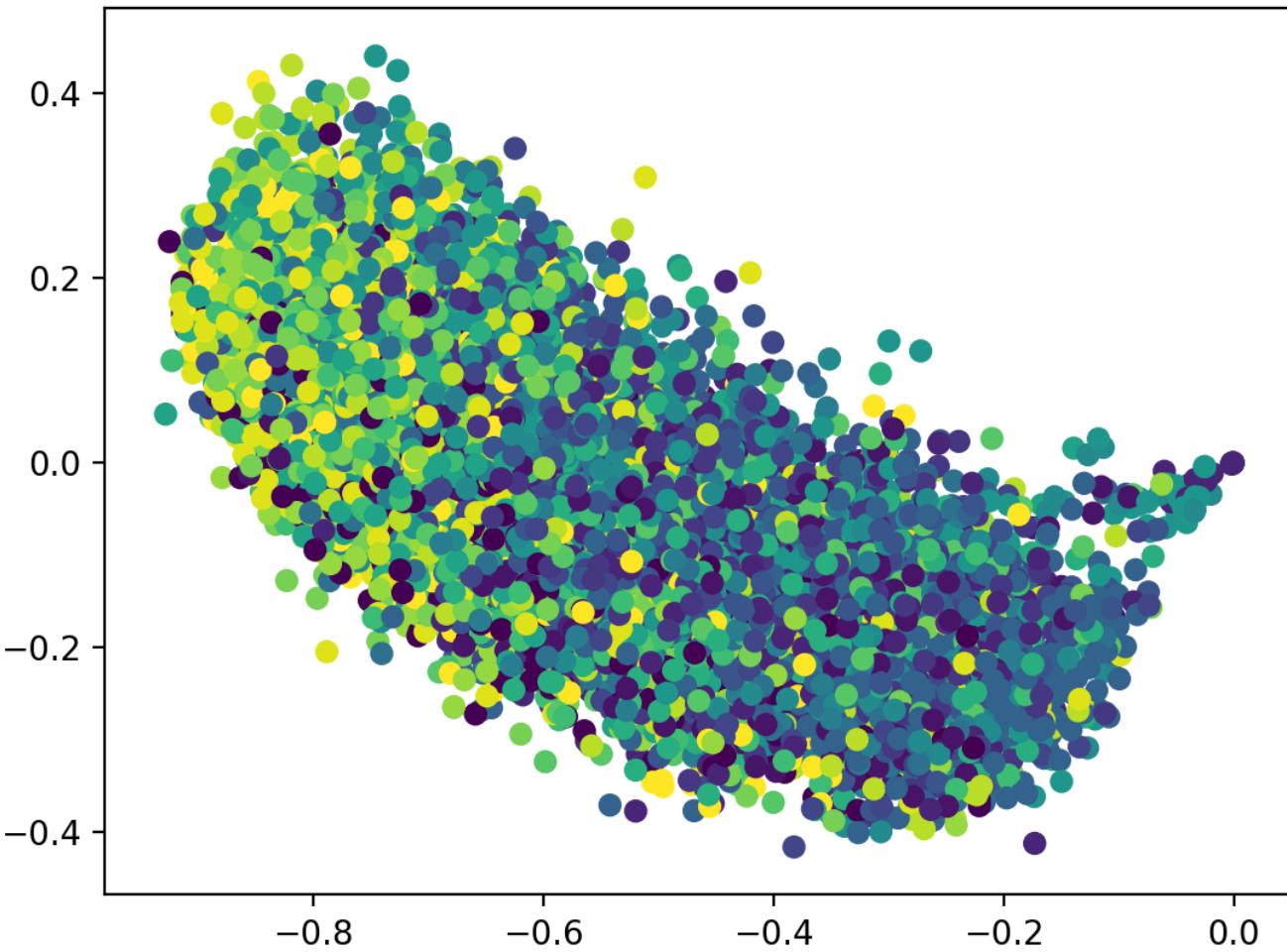
Part(a) =

The code for the program is included  
with the assignment submission.

Code file Name =

LSI-PCA-orthogonal-iteration.py

By running this program following plot  
for 20newsgroup Dataset is obtained =



Different intensities of green color  
represent 20 different clusters of  
newsgroups.

Part (b) =

The code file is added with the assignment.

This prints top ten terms with highest value of ( $\theta_{ik}$ ) for each cluster.

Higher the value means the word is most commonly used in the articles of that newsgroup cluster.